

## Online Convex Optimization (OMD, Duality, Regret)

*Lecturer: Kris Kitani**Scribes: Haowen Shi, Fan Jia*

## 1 Review

Previously we studied Follow The Leader (FTL) and Follow The Regularized Leader (FTRL) algorithms in the Online Convex Optimization (OCO) [2] category.

### 1.1 Follow The Leader with Quadratic Loss

---

**Algorithm 1** FTL with Quadratic Loss
 

---

```

1: function FOLLOW THE LEADER (QUADRATIC LOSS)
2:    $f^{(0)} \leftarrow 0$ 
3:   for  $t = 1, 2, \dots, T$  do
4:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in S} \sum_{i=1}^{t-1} f^{(i)}(\mathbf{w})$ 
5:     RECEIVE ( $f^{(t)} = \frac{1}{2} \|\mathbf{w} - \mathbf{z}^{(t)}\|_2^2$ ) ▷ Receive loss function
6:   end for
7: end function
  
```

---

We proved the regret bound is:

$$\text{Regret} \leq 4L^2 (\log(T) + 1)$$

This is no-regret since the average regret grows sub-linearly in  $T$ . However, this only holds for quadratic loss functions, not for all convex functions.

### 1.2 Follow The Leader with Linear Loss

---

**Algorithm 2** Follow the leader algorithm with linear loss
 

---

```

1: function FOLLOW THE LEADER (LINEAR LOSS)
2:   for  $t = 1, 2, \dots, T$  do
3:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in S} \left( \sum_{i=1}^{t-1} \mathbf{z}^{(i)} \right) \cdot \mathbf{w}$ 
4:     RECEIVE ( $f^{(t)}(\mathbf{w}) = \mathbf{z}^{(t)} \cdot \mathbf{w}$ ) ▷ Receive loss function
5:   end for
6: end function
  
```

---

We observed that the regret of FTL with linear loss relative to an expert  $u$  that chooses  $w = 0$  is:

$$\text{Regret}(u) = T - 1 - 0 = T - 1 = O(T)$$

The regret grows linearly because the  $\mathbf{w}$  value flip-flops back and forth and make the algorithm unstable. To resolve this issue, we introduced regularization.

### 1.3 Follow The Regularized Leader

---

**Algorithm 3** Follow The Regularized Leader

---

```

1: function FTRL(CONVEX SET  $\mathcal{S}$ )
2:   for  $t = 1, 2, \dots, T$  do
3:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in \mathcal{S}} \sum_{i=1}^{t-1} f^{(i)}(\mathbf{w}) + \psi(\mathbf{w})$  ▷ Added  $\psi(\mathbf{w})$  to the prediction rule
4:     RECEIVE ( $f^{(t)} : \mathcal{S} \rightarrow \mathbb{R}$ )
5:   end for
6: end function

```

---

We proved the regret bound for FTRL is:

$$\begin{aligned}
R(\mathbf{u}) &\triangleq \sum_{t=1}^T \left[ f^{(t)}(\mathbf{w}^{(t)}) - f^{(t)}(\mathbf{u}) \right] \\
&\leq \left[ \psi(\mathbf{u}) - \psi(\mathbf{w}^{(1)}) \right] + \sum_{t=1}^T \left[ f^{(t)}(\mathbf{w}^{(t)}) - f^{(t)}(\mathbf{w}^{(t+1)}) \right]
\end{aligned}$$

---

**Algorithm 4** FTRL with Linear Loss And Quadratic Regulation

---

```

1: function FTRL(CONVEX SET  $\mathcal{S}$ )
2:   for  $t = 1, 2, \dots, T$  do
3:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in \mathcal{S}} \sum_{i=1}^{t-1} \left( \frac{1}{2\eta} \|\mathbf{w}\|_2^2 + \mathbf{w} \cdot \sum_{i=1}^{t-1} z^{(i)} \right)$  ▷ Prediction rule
4:     RECEIVE ( $f^{(t)} : \mathbf{w} \cdot z^{(t)}$ )
5:   end for
6: end function

```

---

We proved the regret bound for FTRL with a Euclidean regularizer and linear loss function is:

$$R^{(T)}(\mathbf{u}) \leq BL\sqrt{2T}$$

Where  $L = \max_{\mathcal{Z}} \|\mathbf{z}\|_2$  and  $B = \max_{\mathbf{u} \in \mathcal{S}} \|\mathbf{u}\|_2$ .

*Online Mirror Descent* (OMD) can be considered as a special case of FTRL with a linear loss and a convex regularizer. In this lecture, we will focus on OMD.

## 2 Online Mirror Descent

*Online Mirror Descent* (OMD) [1] can be considered as a special case of FTRL. Fig. 1 shows the relationship between FTL, FTRL and OMD. FTRL adds a regularizer by forcing the zero-step loss to be a non-zero regularization function  $\psi(\mathbf{w})$ . In particular, if a FTRL problem has a convex regularization  $\psi(\mathbf{w})$  and a linear loss function, it is called OMD.

With a linear loss, we can write the loss function as a dot product (linear combination) of the loss function parameter  $\mathbf{z}$  and primal space parameter  $\mathbf{w}$ , as shown in Algorithm 5.

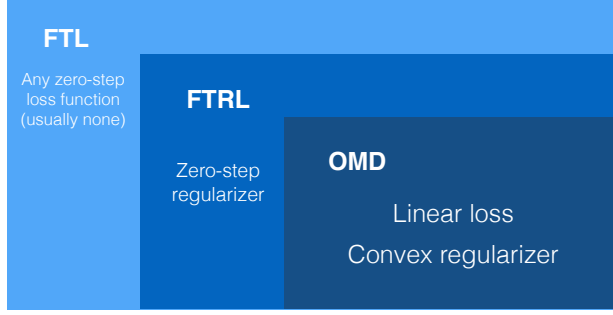


Figure 1: OMD can be seen as FTRL with convex regularization and linear loss

---

**Algorithm 5** Follow The Regularized Leader - Linear Loss

---

```

1: function FTRL-LINLOSS(CONVEX SET  $\mathcal{S}$ )
2:   for  $t = 1, 2, \dots, T$  do
3:      $\mathbf{w}^{(t)} = \arg \min_{\mathbf{w} \in \mathcal{S}} \sum_{i=1}^{t-1} \langle \mathbf{w}, \mathbf{z}^{(i)} \rangle + \psi(\mathbf{w})$             $\triangleright$  Added  $\psi(\mathbf{w})$  to the prediction rule
4:     RECEIVE ( $f^{(t)} : \mathcal{S} \rightarrow \mathbb{R}$ )
5:   end for
6: end function

```

---

We generalize FTRL-LinLoss to derive OMD. First, we define some notations:

Denote sum of  $\mathbf{z}$ :

$$\mathbf{z}^{(1:t)} = \sum_{i=1}^t \mathbf{z}^{(i)} \quad (1)$$

Define parameter of the dual space  $\boldsymbol{\theta}$ :

$$\boldsymbol{\theta}^{(t+1)} \triangleq -\mathbf{z}^{(1:t)} \quad (2)$$

So we have the following update in each iteration for  $\boldsymbol{\theta}$ :

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \mathbf{z}^{(t)} \quad (3)$$

Then, because the loss function  $f^{(i)}$  is linear, we can generalize the FTRL prediction step as follows:

$$\begin{aligned}
\mathbf{w}^{(t+1)} &= \arg \min_{\mathbf{w}} \sum_{i=1}^t f^{(i)}(\mathbf{w}) + \psi(\mathbf{w}) && \text{FTRL prediction step} \\
&= \arg \min_{\mathbf{w}} \langle \mathbf{w}, \mathbf{z}^{(1:t)} \rangle + \psi(\mathbf{w}) && \text{Linear loss function} \\
&= \arg \max_{\mathbf{w}} \langle \mathbf{w}, -\mathbf{z}^{(1:t)} \rangle - \psi(\mathbf{w}) && \text{Convert to maximization} \\
&= \arg \max_{\mathbf{w}} \langle \mathbf{w}, \boldsymbol{\theta}^{(t+1)} \rangle - \psi(\mathbf{w}) && \text{Substitute (2)} \\
&= g(\boldsymbol{\theta}^{(t+1)}) && g: \text{mirror/linking function}
\end{aligned}$$

The function  $g : \boldsymbol{\theta} \rightarrow \boldsymbol{w}$  is called the mirror/linking function because it takes parameters from the dual space ( $\boldsymbol{\theta}$ ) and maps to the primal space ( $\boldsymbol{w}$ ). This function is where the “mirror” name in the OMD came from. With the definitions introduced above, we can change some notations (without changing the algorithm itself) and re-write Algorithm 5 as:

---

**Algorithm 6** Online Mirror Descent

---

```

1: function OMD(CONVEX SET  $\mathcal{S}$ ,  $g : \mathbb{R}^D \rightarrow \mathcal{S}$ )
2:   for  $t = 1, 2, \dots, T$  do
3:     RECEIVE ( $f^{(t)} : \mathcal{S} \rightarrow \mathbb{R}$ )
4:      $\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \eta \boldsymbol{z}^{(t)}$ ,  $\boldsymbol{z} \in \partial f^{(t)}(\boldsymbol{w}^{(t)})$  ▷ Dual parameter update
5:      $\boldsymbol{w}^{(t+1)} = g(\boldsymbol{\theta}^{(t+1)})$  ▷ Primal parameter update (or mirror projection)
6:   end for
7: end function

```

---

The regret bound of OMD (discussed in section 4) shares the same structure with FTRL because at the core they are the same algorithm. With OMD we are doing optimization in the dual space and map back to the primal space to get the updated parameter that we are optimizing (as illustrated in Fig. 2). The benefit of introducing dual space is primarily for theoretical purposes. By looking at the optimization in a different space, we can better understand optimization in the primal space and make more informed changes and improvements.

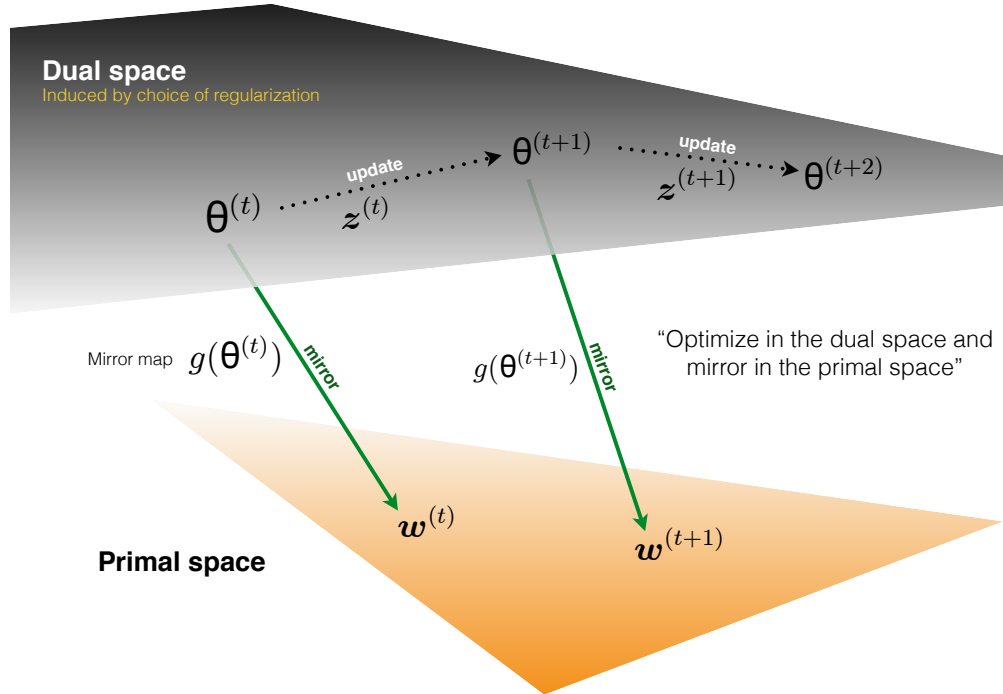


Figure 2: Illustration of online mirror descent in the dual space

A different regularization function  $\psi$  leads to a different mirror function  $g$ . The regularization should be chosen such that the mirroring can take better advantage of the “geometry” of the problem.

## 3 Duality

In this section we introduce mathematical tools that help us analyze OMD.

### 3.1 Convex Conjugate

A *convex conjugate* function is also called a Fenchel dual / a Fenchel conjugate / a Fenchel transform [4]. Here we focus on conjugate functions assuming a smooth convex function. The generalized version of this process for non-smooth and non-convex functions is called a Legendre transformation. The definition for a convex conjugate function is:

$$\psi^*(\theta) = \max_w (\langle \theta, w \rangle - \psi(w)) \quad (4)$$

For better understanding, the geometric interpretation of convex conjugate is included in the appendix section 5.1. Given (4), we can derive some useful properties used to analyze the regret of OMD:

Derivative of convex conjugate:

$$\nabla_{\theta} \psi^*(\theta) = \frac{\partial \psi^*(\theta)}{\partial \theta} = \mathbf{w}^* = \arg \max_w (\langle \theta, \mathbf{w} \rangle - \psi(\mathbf{w})) \quad (5)$$

$$\nabla_{\mathbf{w}} \psi(\mathbf{w}) = \left. \frac{\partial \psi(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\mathbf{w}^*} = \theta \quad (6)$$

Fenchel-Young inequality:

$$\psi^*(\theta) \geq (\langle \mathbf{w}, \theta \rangle - \psi(\mathbf{w})) \quad (7)$$

The proof for (5) is in section 5.2. The proof for (7) is trivial because any  $w$  other than the optimal  $w^*$  will produce less distance  $\psi$  than maximum  $\psi^*$  according to the maximizer definition (4).

### 3.2 Bregman Divergence

Bregman divergence is the approximation error of the first order approximation of  $\psi(\mathbf{w})$  at  $\mathbf{u}$ .

$$D_{\psi}(\mathbf{w}||\mathbf{u}) = \psi(\mathbf{w}) - \psi(\mathbf{u}) - \nabla \psi(\mathbf{u})^{\top} (\mathbf{w} - \mathbf{u})$$

The geometric interpretation of Bregman divergence is included in the appendix section 5.3.

## 4 OMD Regret Analysis

We analysis the regret bound of the OMD algorithm in this section. The regret is:

$$R(\mathbf{u}) = \sum_{t=1}^T \mathbf{w}^{(t)} \cdot \mathbf{z}^{(t)} - \mathbf{u} \cdot \mathbf{z}^{(t)}$$

where  $\mathbf{z}$  is the derivative of the loss function. The following regret bound holds for any OMD algorithm:

$$R(\mathbf{u}) \leq \psi(\mathbf{u}) - \psi(\mathbf{w}^{(1)}) + \sum_{t=1}^T D_{\psi^*}(-\mathbf{z}^{(1:t)} || -\mathbf{z}^{(1:t-1)})$$

The regret is upper bounded by the sum of the difference in the regularization function plus the sum of all Bregman divergence under the convex conjugate of the regularization function. The proof is as follows:

*Proof.* The loss of an arbitrary vector  $\mathbf{u}$ :

$$\psi(\mathbf{u}) + \sum_{t=1}^T \mathbf{u} \cdot \mathbf{z}^{(t)} = \psi(\mathbf{u}) - \mathbf{u} \cdot \boldsymbol{\theta}^{T+1}$$

where  $\boldsymbol{\theta}^{T+1}$  is the sum of gradients till to time step  $T$ .

Recall Fenchel-Young Inequality:

$$\psi^*(\boldsymbol{\theta}) \geq (\langle \mathbf{w}, \boldsymbol{\theta} \rangle - \psi(\mathbf{w}))$$

Apply Fenchel-Young Inequality to the right hand side:

$$\psi(\mathbf{u}) - \mathbf{u} \cdot \boldsymbol{\theta}^{T+1} \geq -\psi^*(\boldsymbol{\theta}^{T+1})$$

Then we expand the right hand side with telescoping. A telescoping series is a series whose partial sums eventually only have a fixed number of terms after cancellation.

Convert to telescoping series by adding and subtracting the complex conjugate at time step 1:

$$-\psi^*(\boldsymbol{\theta}^{T+1}) = -\psi^*(\boldsymbol{\theta}^{T+1}) - \psi^*(\boldsymbol{\theta}^{(1)}) + \psi^*(\boldsymbol{\theta}^{(1)})$$

Sum over shifted time steps:

$$\begin{aligned} -\psi^*(\boldsymbol{\theta}^{T+1}) &= -\psi^*(\boldsymbol{\theta}^{T+1}) - \psi^*(\boldsymbol{\theta}^{(T)}) + \psi^*(\boldsymbol{\theta}^{(T)}) - \dots - \psi^*(\boldsymbol{\theta}^{(1)}) + \psi^*(\boldsymbol{\theta}^{(1)}) \\ &= -\psi^*(\boldsymbol{\theta}^{(1)}) - \sum_{t=1}^T (\psi^*(\boldsymbol{\theta}^{(t+1)}) - \psi^*(\boldsymbol{\theta}^{(t)})) \end{aligned}$$

Plug in definition of Bregman divergence:

$$-\psi^*(\boldsymbol{\theta}^{T+1}) = -\psi^*(\boldsymbol{\theta}^{(1)}) - \sum_{t=1}^T \left( \nabla \psi^*(\boldsymbol{\theta}^{(t)}) \cdot (\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}) + D_{\psi^*}(\boldsymbol{\theta}^{(t+1)} || \boldsymbol{\theta}^{(t)}) \right)$$

The first term can be written as  $\psi(\mathbf{w}^{(1)})$  since:

$$\begin{aligned} \psi^*(\boldsymbol{\theta}^{(1)}) &= \psi^*(\mathbf{z}^{(0)}) && \text{from definition } \boldsymbol{\theta}^{(t+1)} \triangleq -\mathbf{z}^{(1:t)} \\ &= \psi^*(\mathbf{0}) && \text{From definition of conjugate} \\ &= \max_{\mathbf{w}} (\langle \mathbf{w}, \mathbf{0} \rangle - \psi(\mathbf{w})) && \text{Compute dot product} \\ &= \max_{\mathbf{w}} (\mathbf{0} - \psi(\mathbf{w})) && \text{Sign flip max-min conversion} \\ &= -\min_{\mathbf{w}} \psi(\mathbf{w}) && \text{Plug in minimizer of cumulative loss at time 1} \\ &= -\psi(\mathbf{w}^{(1)}) && \text{Minimizer is the primal iterate at step 1} \\ &&& \text{('the one-step look ahead cheater')} \end{aligned}$$

Therefore, we have:

$$-\psi^*(\boldsymbol{\theta}^{(T+1)}) = \psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \left( \nabla \psi^*(\boldsymbol{\theta}^{(t)}) \cdot (\boldsymbol{\theta}^{(t+1)} - \boldsymbol{\theta}^{(t)}) + D_{\psi^*}(\boldsymbol{\theta}^{(t+1)} \| \boldsymbol{\theta}^{(t)}) \right)$$

Plug in definition of dual parameter  $\boldsymbol{\theta}^{(t+1)} \triangleq -\mathbf{z}^{(1:t)}$ :

$$-\psi^*(\boldsymbol{\theta}^{(T+1)}) = \psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \left( \nabla \psi^*(-\mathbf{z}^{(1:t-1)}) \cdot (-\mathbf{z}^{(1:t)} + \mathbf{z}^{(1:t-1)}) + D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \right)$$

Since  $\nabla_{\boldsymbol{\theta}} \psi^*(\boldsymbol{\theta}) = \frac{\partial \psi^*(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \mathbf{w}^*$  and  $-\mathbf{z}^{(1:t)} + \mathbf{z}^{(1:t-1)} = -\mathbf{z}^{(t)}$ :

$$-\psi^*(\boldsymbol{\theta}^{(T+1)}) = \psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \left( \langle \mathbf{w}^{(t)}, -\mathbf{z}^{(t)} \rangle + D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \right)$$

Move negative out of sum, we can get:

$$-\psi^*(\boldsymbol{\theta}^{(T+1)}) = \psi(\mathbf{w}^{(1)}) + \sum_{t=1}^T \left( \langle \mathbf{w}^{(t)}, \mathbf{z}^{(t)} \rangle - D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \right)$$

Therefore:

$$\psi^*(-\mathbf{z}^{(1:T)}) = -\psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \left( \langle \mathbf{w}^{(t)}, \mathbf{z}^{(t)} \rangle - D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \right)$$

Recall the Fenchel-Young Inequality, we can have:

$$\begin{aligned} \langle \mathbf{u}, -\mathbf{z}^{(1:T)} \rangle - \psi(\mathbf{u}) &\leq \psi^*(-\mathbf{z}^{(1:T)}) \\ \langle \mathbf{u}, -\mathbf{z}^{(1:T)} \rangle - \psi(\mathbf{u}) &\leq -\psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \left( \langle \mathbf{w}^{(t)}, \mathbf{z}^{(t)} \rangle - D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \right) \\ \langle \mathbf{u}, -\mathbf{z}^{(1:T)} \rangle - \psi(\mathbf{u}) &\leq -\psi(\mathbf{w}^{(1)}) - \sum_{t=1}^T \langle \mathbf{w}^{(t)}, \mathbf{z}^{(t)} \rangle + \sum_{t=1}^T D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \end{aligned}$$

Rearranging this equation, we can have:

$$\sum_{t=1}^T \langle \mathbf{w}^{(t)}, \mathbf{z}^{(t)} \rangle - \langle \mathbf{u}, \mathbf{z}^{(1:T)} \rangle \leq \psi(\mathbf{u}) - \psi(\mathbf{w}^{(1)}) + \sum_{t=1}^T D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)})$$

Therefore, our general OMD regret bound is:

$$\begin{aligned} R(\mathbf{u}) &= \sum_{t=1}^T \mathbf{w}^{(t)} \cdot \mathbf{z}^{(t)} - \mathbf{u} \cdot \mathbf{z}^{(t)} \\ &\leq \psi(\mathbf{u}) - \psi(\mathbf{w}^{(1)}) + \sum_{t=1}^T D_{\psi^*}(-\mathbf{z}^{(1:t)} \| -\mathbf{z}^{(1:t-1)}) \end{aligned}$$

□

## References

- [1] G. Lan, A. Nemirovski, and A. Shapiro. Validation analysis of mirror descent stochastic approximation method. *Mathematical programming*, 134(2):425–458, 2012.
- [2] S. Shalev-Shwartz et al. Online learning and online convex optimization. *Foundations and trends in Machine Learning*, 4(2):107–194, 2011.
- [3] Wikipedia contributors. Bregman divergence — Wikipedia, the free encyclopedia, 2021. [Online; accessed 03-Mar-2021].
- [4] Wikipedia contributors. Convex conjugate — Wikipedia, the free encyclopedia, 2021. [Online; accessed 03-Mar-2021].



## 5 Appendix (Covered in Lecture)

### 5.1 Geometric interpretation of convex conjugate

#### 5.1.1 Intercept-slope parameterization

There is more than one way to parameterize a function. One method that we are familiar with the most is the function-value parameterization, where a single output value is provided for each input. Another way we can parameterize a function is called a intercept-slope parameterization.

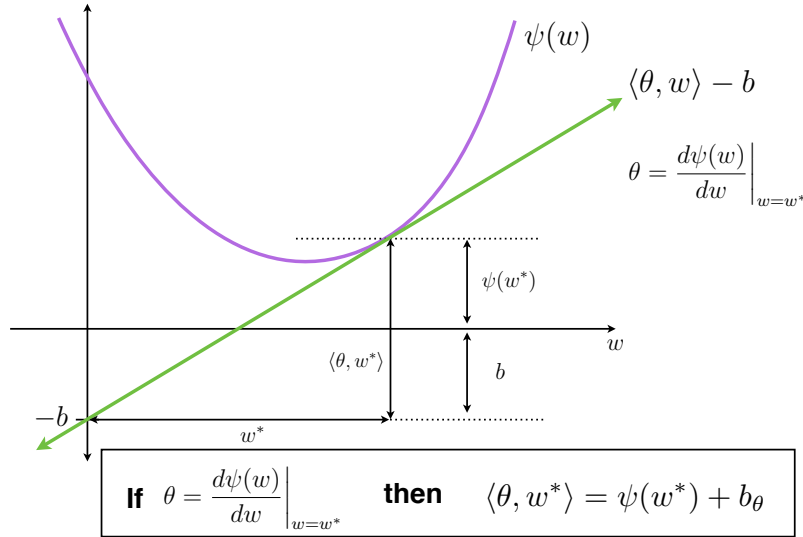


Figure 3: Geometry of tangents

As shown in Fig. 3, we have a function  $\psi(w)$ , and we define  $\theta$  as the derivative evaluated at  $w = w^*$ . Then, the tangent line can be written as  $\langle \theta, w \rangle - b$ . So, the vertical distance (“rise”) between  $-b$  and the tangent point is  $\langle \theta, w \rangle$ . Additionally, the vertical distance can also be written as  $\psi(w^*) + b$ , because of the geometry shown in the figure. So, we have:

$$\langle \theta, w^* \rangle = \psi(w^*) + b(\theta) \quad (8)$$

Where  $b(\theta)$  depends on the  $\theta$  at the chosen  $w^*$ . Rearrange (8), given  $\theta$ , we have:

$$-b(\theta) = -\langle \theta, w^*(\theta) \rangle + \psi(w^*(\theta)) \quad (9)$$

Here we have a one-to-one mapping from  $\theta$  to  $w^*$  because we assume  $\psi$  is convex and smooth. Otherwise, we cannot write  $w^*$  as a function of  $\theta$ .

#### 5.1.2 Conjugate

Now we discuss what is a conjugate function. Consider a slope line passing through the origin:  $\langle \theta, w \rangle$ , and the function  $\psi(w)$ . As illustrated in Fig. 4, we want to look for the maximum distance between the two functions, i.e.  $w_\theta^* = \arg \max_w (\langle \theta, w \rangle - \psi(w))$ . One observation we can make is

that at  $w_\theta^*$ , the slope of the tangent line is same as  $\theta$ . If we draw the tangent line (dotted), the intercept of the tangent line is actually equal to the maximum distance between two functions at  $w_\theta^*$ .

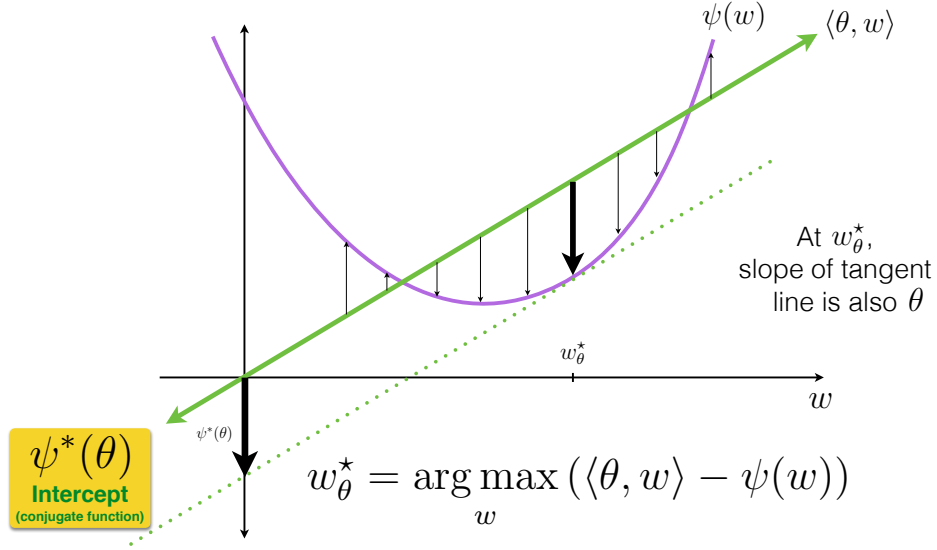


Figure 4: Geometry of the conjugate function  $\psi^*(\theta)$

We call this maximum distance the “intercept”  $\psi^*(\theta)$  and for any slope  $\theta$ , we have:

$$\psi^*(\theta) = \max_w (\langle \theta, w \rangle - \psi(w)) \quad (10)$$

This is the geometric interpretation of the conjugate function (4). To reiterate, the convex conjugate at  $\theta$  is the maximum distance between function  $\psi$  and the slope function  $\langle \theta, w \rangle$  passing through the origin.

## 5.2 Proofs for properties of convex conjugate

Proof for (5):

*Proof.*

$$\begin{aligned} \psi^*(\theta) &= \max_w (\langle \theta, w \rangle - \psi(w)) && \text{Definition (10)} \\ \nabla_\theta \psi^*(\theta) &= \nabla_\theta \max_w (\langle \theta, w \rangle - \psi(w)) && \text{Take derivative} \\ &= \nabla_\theta (\langle \theta, w^* \rangle - \psi(w^*)) && \text{In terms of optimal } w^* \\ &= w^* && \text{Compute partial derivative} \end{aligned}$$

□

### 5.3 Geometric interpretation of the Bregman divergence

Fig. 5 shows the geometric interpretation of Bregman divergence [3].  $\psi(\mathbf{w}) - \psi(\mathbf{u})$  is difference of function values.  $\nabla\psi(\mathbf{u})^\top(\mathbf{w} - \mathbf{u})$  is the first order approximation of  $\psi(\mathbf{w})$ .  $D_\psi(\mathbf{w}||\mathbf{u})$  is the Bregman divergence. Geometrically, it is the “distance” between the point on the function  $\psi$  (regularization function in the case of OMD) and its tangent at  $\mathbf{u}$  when evaluated at  $\mathbf{w}$ . It represents the first order approximation error evaluated at  $\mathbf{w}$  with linearization point  $\mathbf{u}$ .

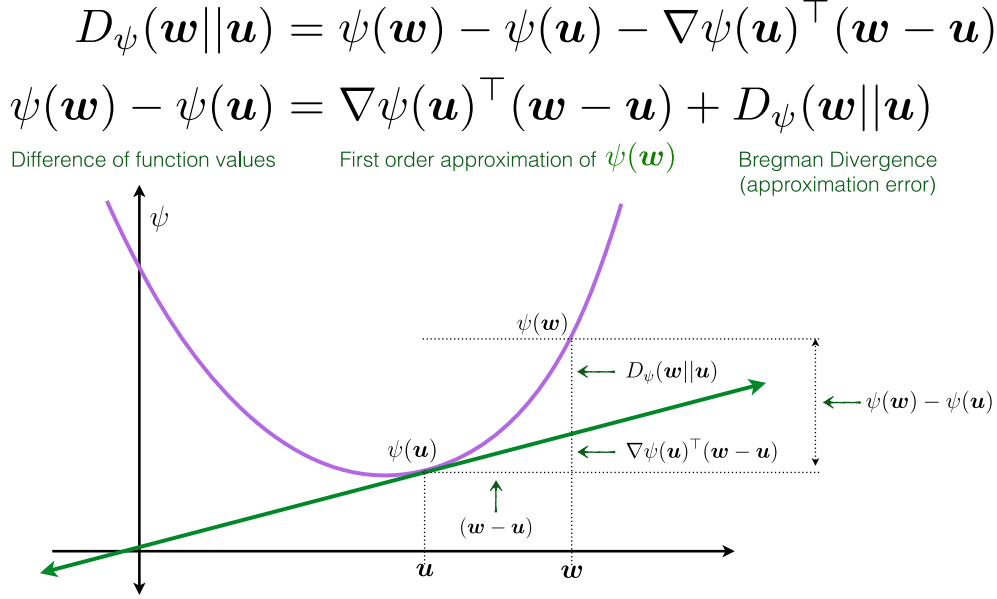


Figure 5: Geometric interpretation of Bregman divergence